

基于 IIIF 和语义知识图谱的印章资源整合与知识发现研究*

■ 张永娟^{1,2} 刘炜³ 于建荣² 陈涛^{3,4}¹ 上海大学图书情报档案系 上海 200444² 中国科学院上海营养与健康研究所/上海生命科学信息中心 上海 200031³ 上海图书馆/上海科学技术情报研究所 上海 200031 ⁴ 南京大学信息管理学院 南京 210023

摘要: [目的/意义] 数字人文研究的图像资源中蕴含大量信息但利用率极低,不能在异构数据库和不同的应用程序中得到有效的共享与重用,国际图像互操作框架打破了图像资源交换和共享的障碍。[方法/过程] 研究结合国际图像互操作框架和语义知识图谱(关联数据技术)进行图像资源的整合、共享与知识发现,对资源之间的关系进行揭示和知识推理,并通过 CNNs 算法对图像特征的提取与识别实现基于图像特征的语义检索辅助知识发现。[结果/结论] 提出一套数字人文图像资源整合与知识发现解决方案,并以印章图像资源为应用对象构建“印章知识中心”对以上解决方案的可行性和实践性进行实证检验。

关键词: 数字人文 图像资源整合 IIIF 关联数据 知识图谱 知识发现

分类号: G203

DOI: 10.13266/j.issn.0252-3116.2020.07.015

1 引言

数字人文(Digital Humanities)是现代信息技术应用于传统人文研究而形成的一个新型的跨学科研究领域,近年来获得广泛关注,在国内外掀起了研究和应用的热潮^[1]。图像作为数字人文领域一种传递信息、知识和思想的非文本视觉媒介,具体表现形式十分多样,包括绘画、照片、草图、手稿、印章等,图像包含了深刻的文化内涵、复杂的时空场景和较为抽象的思想寓意。目前在特定领域图像注释模型有利于用户对图像的理解,但在不同领域图像资源的共享、复用、整合及知识发现等方面还存在壁垒和障碍。图像仍然被禁锢在数据库中,图片无法共享和复用,国际图像互操作框架(International Image Interoperability Framework, IIIF)在国内应用仍处于探索阶段,国际互操作标准、关联数据、知识图谱等相关技术在图像资源方面的应用也较少。图像的共享、复用、整合与知识发现成为国内数字人文领域亟待解决的重要问题之一。

2 研究现状

2.1 IIIF 在图像资源语义互操作方面的研究进展

在图像数据互操作方面,国外有着较为丰富的理论研究和实践应用经验。IIIF^[2]于2015年由欧洲和美国的图书馆等29个非营利图像资源存储机构共同成立,对以图像为载体的书籍、地图、卷轴、手稿、乐谱、档案等在线资源进行统一展示和共享。国际上主要的文化遗产研究机构都采用了IIIF对其图像进行管理和共享。IIIF解决了文化资源数字化图像难以被发现、再利用、引用、交换、比较分析等问题,为确保全球图像存储的互操作性和可获取性提供了国际化通用标准。

在数据注释方面也存在成熟的国际标准,开放协同标注(Open Annotation Collaboration, OAC)^[3]是最早提出的促进标注的规范化、共享和复用的国际标准。W3C的开放注释数据模型(Open Annotation Data Model, OADM)在OAC的基础上引入了关联数据技术,其作为数据注释的国际互操作框架,以众包的方式提供在线的语义标注,允许数字人文研究学者添加更多的

* 本文系中国博士后科学基金项目“数字人文中印章图像的深度学习与语义检索”(项目编号:2018M2058)研究成果之一。

作者简介: 张永娟(ORCID:0000-0002-7776-547X),副研究馆员,硕士;刘炜(ORCID:0000-0003-2663-7539),研究员,博士;于建荣(ORCID:0000-0003-0843-1764),研究馆员,硕士;陈涛(ORCID:0000-0002-6609-4914),工程师,博士后,通讯作者;E-mail: tchen@lib-net.sh.cn。

收稿日期:2019-08-29 **修回日期:**2019-10-26 **本文起止页码:**127-135 **本文责任编辑:**徐健

关联,丰富其内容,并基于关联数据这个将数据开放到互联网的方法创建资源、注释之间的关联,实现平台间资源和注释共享与开放。

学者曾蕾在相关会议上曾多次提出使用开放协同标注和 IIIF 实现图像语义深度标引的建议,对这些国际标准在国内的应用起到引导和推动作用^[4]。上海图书馆、武汉大学、北京大学、复旦大学以及上海慧游文化传播有限公司等机构都在积极探索基于 IIIF 的图片管理与共享方案。目前虽尚未见基于 IIIF 和关联数据的图像资源的大规模的应用报道,但随着技术的突破,IIIF 在数字人文领域图像资源的应用可能很快进入一个爆发点。

2.2 语义知识图谱(关联数据)在图像资源整合与知识发现方面的应用

知识图谱是利用计算机存储、管理和呈现概念及概念间关系的一种技术,可分为基于 RDF 存储的语义知识图谱(即关联数据)和基于图数据库的广义知识图谱。语义知识图谱(关联数据)侧重于知识的发布和链接,广义知识图谱更侧重于知识的挖掘和计算,关联数据是谷歌知识图谱的延续和发展,广义知识图谱研究中丰富的图运算和关联数据的结合将会带来数字人文研究的新时代^[5]。在图情界和数字人文领域,提的较多的是语义知识图谱(关联数据)。

语义知识图谱(关联数据)和 IIIF 互为补充,在图像资源关联整合、共享方面发挥重要作用。伏尔泰书信、达芬奇手稿等档案资源也以关联数据形式和 IIIF 对其图像资源进行语义组织和发布。Linked Canvas^[6] 图像语义注释共享解决方案是 Synaptica 开放注释语义检索系统的(Open Annotation Semantic Indexing System, OASIS)的重要补充,其使用关联数据技术、词表和本体对图像内容进行丰富,基于关联数据平台(LDP)、W3C 的 OADM 数据模型、IIIF 语义互操作框架建立全球文化遗产社区非文本数据与注释数据的组织、关联和共享,使注释能够在不同的硬件和软件平台上共享和重用。

机器学习/神经网络算法等人工智能(AI)技术的进步也推动着语义知识图谱(关联数据)、IIIF、图片语义检索在数字人文领域的发展和进步。“威尼斯时间机器”(Venice Time Machine)^[7]项目是瑞士联邦理工学院(EPFL)数字人文科学实验室利用机器学习算法,将威尼斯多年的历史以动态的数字化形式传承下来,再现这座古城辉煌的共和国时代风貌,揭示整个欧洲大陆当时的社会网络、贸易和知识发展的历史。

国内关于关联数据的研究除了专业期刊上发表的大量研究报告和论文之外,实践应用主要集中在文本数据方面,如上海图书馆推出的家谱知识库、古籍循证平台、名人手稿知识库等一系列关联数据应用平台^[8];曾子明^[9]将关联数据技术应用于敦煌视觉资源关联展示;侯西龙等^[10]将关联数据用于非物质文化遗产知识管理研究中。这些研究也可以看成是知识图谱的应用研究,然而其中大多数应用系统都是使用关联数据技术来进行元数据层面的知识组织和发布,极少使用知识图谱的理念对资源之间的关系进行揭示和知识推理。中国历代人物传记资料库(China Biographical Database Project, CBDB)^[11]借助知识图谱的理念展现了人物之间丰富的亲属及社会关系,形成特有的社会关系网络,并可通过设置推理规则实现人物之间隐性关系的挖掘与呈现。但通过 IIIF、语义知识图谱(关联数据)等实现图像资源的整合、共享和复用的实践案例还鲜有报道。

从总体上来看,国外有将关联数据、IIIF、机器学习算法等 AI 技术同时用于文本和图像资源的整合与知识发现的成熟实践案例,也得到同行的认可,是未来数字人文研究的重要方向之一。国内研究主要集中在关联数据技术对文本信息元数据的描述和研究,对于数字人文领域基于 IIIF、语义知识图谱(关联数据)的图像资源整合和知识发现的应用研究不多,还存在需要突破的技术难点和较大的挑战。因此,本研究尝试运用 IIIF 以及语义知识图谱(关联数据)等语义技术,突破技术难点,建立可行的解决方案,实现图像资源的整合、隐形关系揭示与知识发现。

3 图像资源整合与知识发现解决方案

本研究通过对国外成功的实践案例进行调研和剖析,提出了图像资源整合与知识发现解决方案,并基于印章图像资源构建“印章知识中心”,对解决方案的可实施性进行验证,实现印章图像资源与其他资源的整合,以语义知识图谱的方式实现知识发现。

该解决方案主要涉及基于 IIIF 的图像元数据描述,以及基于语义知识图谱(关联数据)的知识发现两部分,其中基于 IIIF 的图像元数据描述包括图像 API、呈现 API 及图像注释;基于语义知识图谱(关联数据)的知识发现包括 KOS/本体构建、关联数据发布服务、语义索引及语义注释,同时借助深度学习的方法实现图像检索,最终实现资源整合、知识图谱呈现及知识发现服务。解决方案架构如图 1 所示:

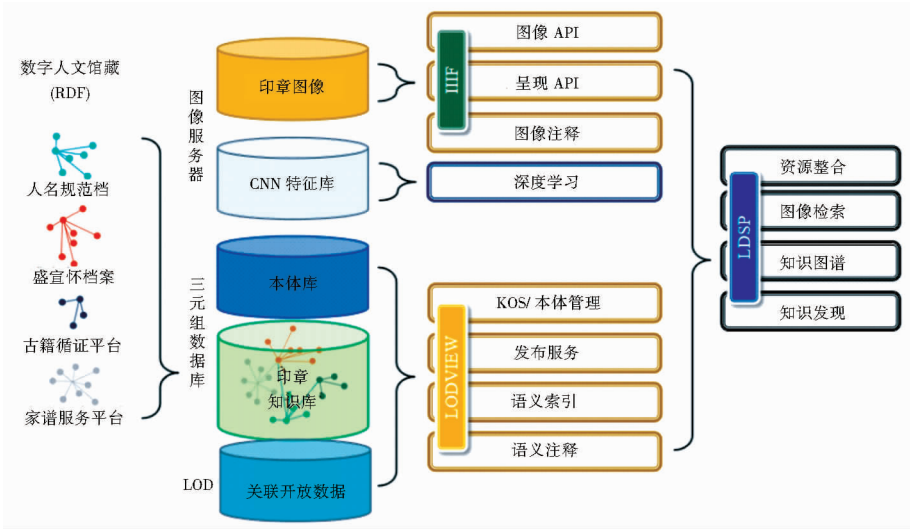


图 1 图像资源整合与知识发现解决方案架构设计

3.1 基于 IIIF 和 OADM 的图像元数据描述

国际图像互操作框架 IIIF 定义了一组通用 API (应用程序编程接口) 规范, 支持图像存储库之间的互操作性。从最基础的系统层, 解决了硬件和操作系统之间不兼容的问题。IIIF 目前拥有 4 个标准 API, 可用于图像元数据的规范、图像的呈现、语义注释与共享、以及语义检索, 也可以根据实际需要扩展新的 API。

“图像 API (Image API)”可作为图像编目过程的一部分, 通过 URI 指定所请求图像的来源、区域、大小、角度、质量、格式等。“呈现 API (Presentation API)”指定一个返回 JSON-LD 结构化文档的 Web 服务, 共同描述数字化对象的结构和布局, 可用于多个图像的比较、使用者在线标注, 实现有出处来源与可控的分享。其结构包括: 整套藏品 (collection)、整个物件 (manifest)、所有张页的顺序 (sequence)、单页 (canvas)、相关关系 (anno)、数字内容 (content), 可为每个图像注释分配一个唯一的 HTTP URI, 实现通过 HTTP 在线访问和注释, 在原有系统上实现图片的递送和有出处来源的共享, 最终可以实现众多资源的重新组合, 实现出版一次复用多次。“搜索 API (search API)^[12]”可支持在单个 IIIF 资源中搜索注释内容。“验证 API (authentication API)”描述了一组用于引导用户完成现有访问控制系统的工作流程。

注释 (annotation)^[13] 是在不同信息之间建立关联的标记行为。W3C 的 OADM 开放注释数据模型提供了一个可扩展的、可互操作的框架用于表达注释, 使得它们可以在平台之间轻松共享, 以最简单的方式满足

最复杂的需求。OADM 为其类和属性定义命名空间, 即使本体发生更改, 命名空间 URI 也将始终保持不变。所有版本的本体都将从特定于版本的 URL 保持可用, 并且命名空间 URI 将提供对最新版本的访问。OADM 开放注释数据模型结合 IIIF 搜索 API 可以实现图像的结构组织与重用 (见图 2), 开放注释数据内容标注在 IIIF 呈现 API 的 canvas 上, 可以标注整个 canvas, 或者部分区域, 区域选择可以是任意形状, 注释支持个人或多人在线协作, 可以众包的形式对用户开放。canvas 作为一个新的交互层和链接在 Web 上, 赋予唯一的 URI, 它允许任何人在任何地方注释任何内容, 无论是网页、电子书、视频、图像、音频流, 还是原始或可视化形式的数据, 其不同服务之间实现标注内容的链接和共享, 并可追溯到它们的起源, 便于搜索和发现。

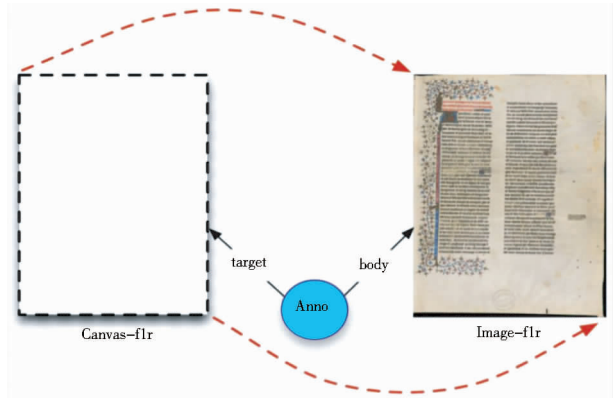


图 2 结合 IIIF 和 OADM 的图像语义注释图解

IIIF 和 OADM 开放注释数据模型等国际开放标准的使用都为进一步实现图像的语义检索和知识图谱的

构建奠定了基础。笔者团队在图像资源整合与知识发现解决方案与“印章知识中心”对 IIIF 图像 API 的使用进行尝试,最终实现印章图像的深度缩放和在线调用,可通过 URI 指定所请求印章图像的来源、区域、大小、角度、质量、格式等;采用 IIIF 呈现 API 描述印章资源的结构和布局,可用于多个印章图像的比较、使用者在线标注,实现有出处来源与可控的分享,OADM 开放注释数据内容标注在 canvas 上,目前采用将整段注释标

注在整个 canvas 上的方式,将 IIIF 呈现 API 和 OADM 注释做了有意义的尝试和探索,在接下来的工作中可采用众包等多人在线协作的形式进行更深入的分区注释,或者基于 OCR 按照字符进行注释;采用 IIIF 搜索 API 实现印章图像的图像检索、注释检索和文本检索;因为印章平台尚未涉及权限的控制问题,所以在平台中没有使用验证 API,后续如果需要可进行扩展,如图 3 所示:

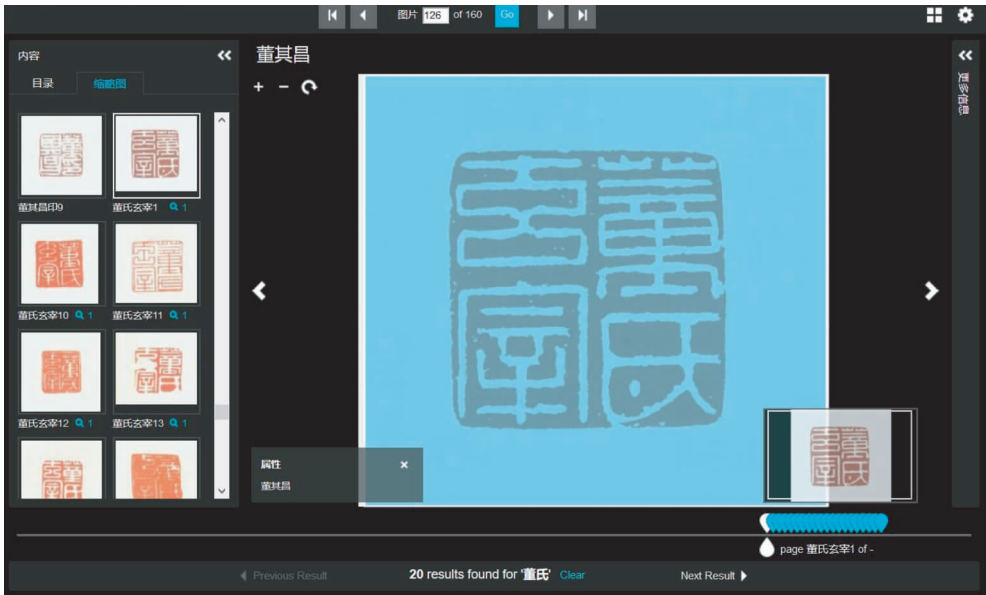


图 3 基于 IIIF APIs 与 OADM 的图像资源呈现

3.2 图像描述内容知识组织与关联数据发布

(1) 图像描述内容的知识组织。解决方案中采用本体用于图像描述内容的知识组织,在本体设计时,尽量复用已有本体的类和属性,依据该原则,印章本体在设计时复用了 sh1、foaf 等这些命名空间下面定义的和属性,详见表 1。另外,印章平台扩展了 owner 用来实现印章平台和上海图书馆人名规范库的资源的关联,同时采用了 sameAs 和 owner 实现其他不同数据集之间的关联。

表 1 用于组织印章内容本体的核心类与属性

名称	类型	描述
sh1:Seal	类	印章类
sh1:sealCharacters	对象属性	印章印文
sh1:owner	对象属性	印章主人
sh1:ownerOfSeal	数据属性	印章主人
foaf:img	对象属性	印章图片

基于本体对印章内容进行知识组织的 RDF 示例如下:

```
<http://data.library.sh.cn/gj/entity/seal/rvwiaqlu1yjzh33> a sh1:Seal;
rdfs:label "董氏玄宰 6"@chs, "董氏玄宰 6"@cht;
foaf:img <http://data.library.sh.cn/gj/resource/img/hy4gev5imtzevnhj>;
foaf:img <http://data.library.sh.cn/gj/resource/img/dkg13vqqspredevqy>;
foaf:img <http://data.library.sh.cn/gj/resource/img/wf1lm863dy8v42a2>;
foaf:img <http://data.library.sh.cn/gj/resource/img/iuyj1ieqbeerc3zu>;
foaf:img <http://data.library.sh.cn/gj/resource/img/pj25proelo13hdq>;
sh1:ownerOfSeal "董其昌";
sh1:sealCharacters "董氏玄宰 6"@chs, "董氏玄宰 6"@cht;
sh1:owner <http://data.library.sh.cn/entity/person/7jfqns5vitt6efhl>.
```

(2) 图像描述内容的关联数据发布。本模型使用 SinoPedia 关联数据发布平台 (SinoPedia Platform, LD-SP)^[14] 将图像内容实现七星标准^[15] 关联数据的发布 (见图 4), LDSP 是本团队的前期研究成果,它不仅可以作为独立的知识库进行资源检索,还可以作为关联数据发布中心 (Linked Data Hub) 来发布多源的关联数据集,并提供相关资源的关联数据发布和内容协商服务。

LDSP 平台提供的关联数据转换服务 (linked data transformation service, LDTS) 可将非结构化、半结构化、结构化的数据转换成关联数据并存储在三元组数据库

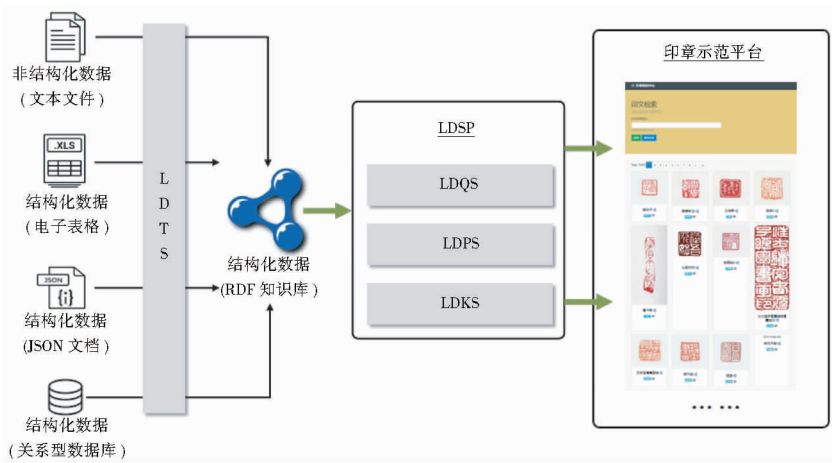


图 4 印章平台调用 LDSP 服务解析

中 (OpenLink Virtuoso); 关联数据查询服务 (linked data query service, LDQS) 为不同的数据集提供 SPARQL 端点, 支持 SPARQL 联合查询; 关联数据发布服务 (linked data publishing service, LDPS), 支持单个 LODVIEW 平台接入多个 SPARQL 端点, 在 SinoPedia 上显示不同站点资源, 并提供关联数据的内容协商服务 (主要有 RDF/XML、JSON_LD、NT、TTL)。LODVIEW 系统是基于 Spring 和 Jena 的 Web 应用程序, 可以与 SPARQL 端点共同按照关联数据的发布标准来发布 RDF 数据, 将图像元数据按照七星数据模型的标准和规范发布出来; 关联数据知识服务 (linked data knowledge service, LDKS) 通过集成 LODLIVE 模块实现对关联的多源数据集进行知识整合和知识图谱展示^[16]。

通过该模型构建的七星关联数据知识库支持多种应用程序的开发和数据调用, 这种基于语义 (概念) 和关联数据知识库的检索方式, 增强了图像资源相关的语义内容。

3.3 图像资源整合、知识图谱与知识发现

(1) 图像资源整合与知识图谱的构建。语义知识图谱 (关联数据) 本质上是一种由知识点相互连接而成的语义网络, 支持搜索引擎进行知识发现、索引以及可视化呈现。本研究借助 LDSP 平台提供关联数据知识服务 (LDKS) (见图 4) 实现知识图谱与知识发现的应用, LDKS 提供的知识图谱和可视化等相关技术, 可将不同知识库 (包括 LOD 中的关联数据集和上海图书馆发布的关联数据集) 中的多源数据集融合。

发布的关联数据知识库可以实现与外部关联数据知识库的关联与融合, 外部的关联数据知识库为图像描述提供更丰富的关联和语义增强, 也为知识图谱可

视化提供丰富的资源。不同主题的数据集可以根据内容关联不同的外源关联数据知识库。

不同数据集资源之间的关联与融合主要是通过 OWL 的 sameAs、seeAlso 等属性, 其中使用较广的 sameAs 用于连接两个实体是相同的本体之间的映射。LDKS 服务将主要数据集的 sameAs 关系抽取到一个中心池 (存放 sameAs 的 graph), 作为中心 “映射” 层的网络基础设施, 统一动态收割 sameAs 属性, 并对有关系的数据集建立双向链接。其他外部数据集与中心库任意数据集建立链接, 将自动获取与其相关的数据集 sameAs 关系。

人物是印章平台的核心要素, 在做印章知识图谱时, 主要以藏印主人为核心关联对象进行不同数据源之间的信息关联与融合, 外部关联数据知识库包括来自关联数据云 (Linked Open Data Cloud, LOD) 的 LOC、VIAF 以及 DBpedia 数据集, 也包括来自上海图书馆开放的人名规范档、古籍知识库、SinoPedia 以及 CBDB 等数字人文研究国家数据基础设施。关联后的数据集经由知识图谱的可视化展示, 可以从多个角度揭示图像背后蕴含的丰富内容, 这些都为图像的语义检索提供了知识来源, 最终实现跨数据集之间资源的知识发现, 帮助用户更好地从关联数据中挖掘、分析隐含知识, 提供多维知识服务, 见图 5。

(2) 基于机器学习的图像检索与知识发现。图像的非文本内容 OCR 很难识别, 印章图像中的内容多为繁体字、古文字, 也给 OCR 的识别带来了难度, 本研究可以尝试采用深度学习的方法来进行图像特征的识别和提取, 最终实现图像检索, 尝试用人工智能的手段辅助知识发现。

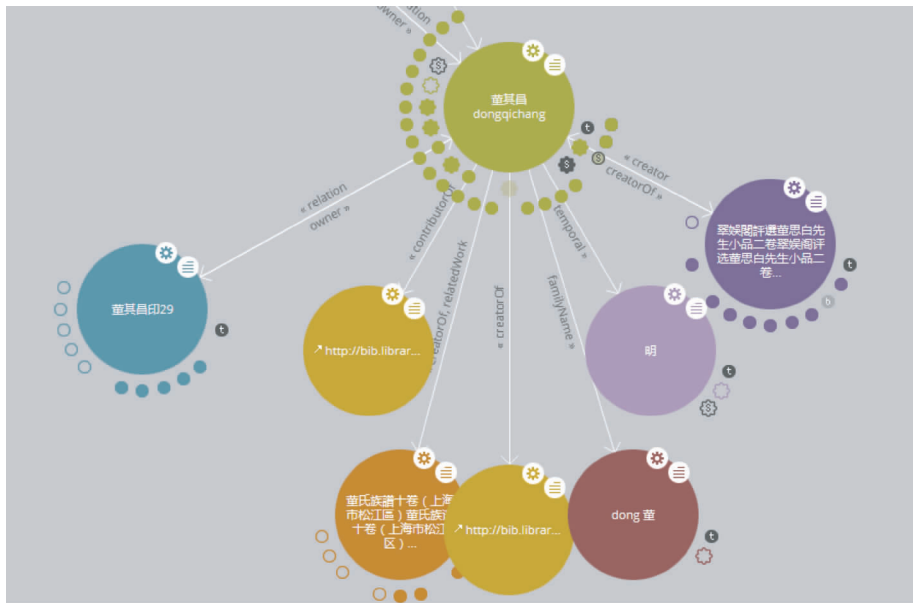


图 5 图像资源整合与知识图谱呈现

卷积神经网络 (Convolutional Neural Networks, 简称 CNNs) 是一种深度的监督学习下的机器学习模型, 具有极强的适应性, 善于挖掘数据局部特征, 提取全局训练特征和分类。它的权值共享结构网络使之更类似于生物神经网络, 在模式识别各个领域都取得了很好的成果。相比于基于 SIFT、HOG 等特征的图像分类方法, 基于 CNNs 的方法有更强的高层语义抽象能力, 同时 CNNs 有着天然的平移不变性和训练可得的一定范围内的尺度不变性, 这些特性也都是图像分类所必须具备的。

本模型采用牛津大学 VGG 组提出的深度 CNNs 模型 VGG16 进行图像特征提取, 它改进了 AlexNet 中的较大卷积核, 采用连续的几个小的卷积核代替, 采用堆积的小卷积核是优于采用大的卷积核, 因为多层非线性层可以增加网络深度来保证学习更复杂的模式, 而且代价还比较小 (参数更少), 在图像分类等任务中取得了不错的效果, 如图 6 所示:

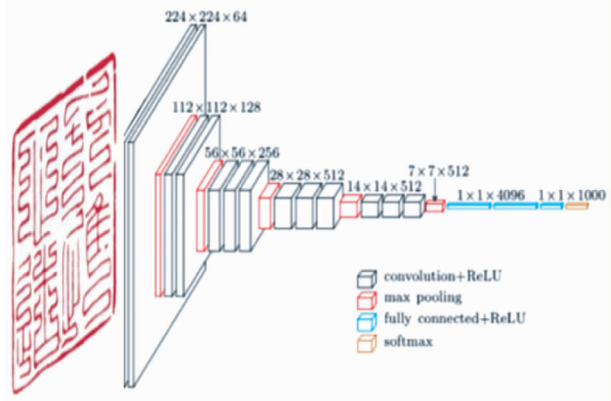


图 6 深度卷积神经网络模型 VGG16 原理

通过使用模型 VGG16 对印章图像特征进行提取, 形成印章图像特征库; 在用户端, 对用户用于检索的图像也通过使用 VGG16 模型, 提取图像特征, 并与已形成印章图像特征库的特征进行对比检索, 完成基于图像特征的图片检索。用户可根据偏差度选择相关检索结果中的图像, 每个图像都关联整合了跨库的相关资源, 点击可以以知识图谱的方式呈现不同的资源之前的关系, 进一步点击可通过知识图谱进行可视化查询, 实现知识发现, 具体实例如图 7 所示:



图 7 基于 VGG16 实现图像检索

4 图像资源整合与知识发现解决方案应用案例分析

印章是一类比较有代表性和研究意义的图像, 除了单独的印章图像, 还广泛存在于大幅的画卷中, 也可

作为画卷的局部进行注释和标记。印章图像中除了印文外,还有很多相关信息,比如印章持有者的人物信息、印章使用的历史背景、印章的变迁等,这些深层次的内容通过语义标注可以极大丰富印章的内容,有利于实现图像资源整合与知识发现。笔者为了采用一些实例数据对图像资源整合与知识发现解决方案进行验证,构建了“印章知识中心”。“印章知识中心”目前共收录 15 053 枚印章,包括爱新觉罗弘历、张大千、董其昌等人,实现了印章知识库与上海图书馆人名规范库、古籍知识库、CBDB 的整合与关联^[17]。

本部分以“印章知识中心”中与董其昌相关的 160 枚印章为应用案例,对模型的可用性和有效性进行分析和验证,用户可登陆“印章知识中心”尝试获取更多数据以及其他人物的知识发现。董其昌(1555-1636),字玄宰,明朝后期大臣,著名书画家。“印章知

识中心”中收录董其昌相关印章 160 枚,目前主要包含印文和印章主人两类信息。董其昌 50 岁以前所用姓名字号印有“董其昌(16 枚)”“董其昌印(43 枚)”“董氏玄宰(20 枚)”;50-59 岁增加“其昌之印(1 枚)”及“董玄宰(12 枚)”,60-69 岁增加“其昌”“玄宰(21 枚)”“思白”;69 岁开始有“昌”字印,80 岁后用“思翁(2 枚)”印。

“印章知识中心”通过 SPARQL 联合查询可以实现关联的知识库中的隐含关系,发现背后的知识,如图 8 通过 SPARQL 联合查询,以及扩展的 owner 实现印章平台与上海图书馆人名规范库的关联,通过存放 sameAs 的 graph 实现上海图书馆人名规范库与 CBDB 知识库的关联,进而发现“董其昌的妻子是明代成岫”这一隐含知识。

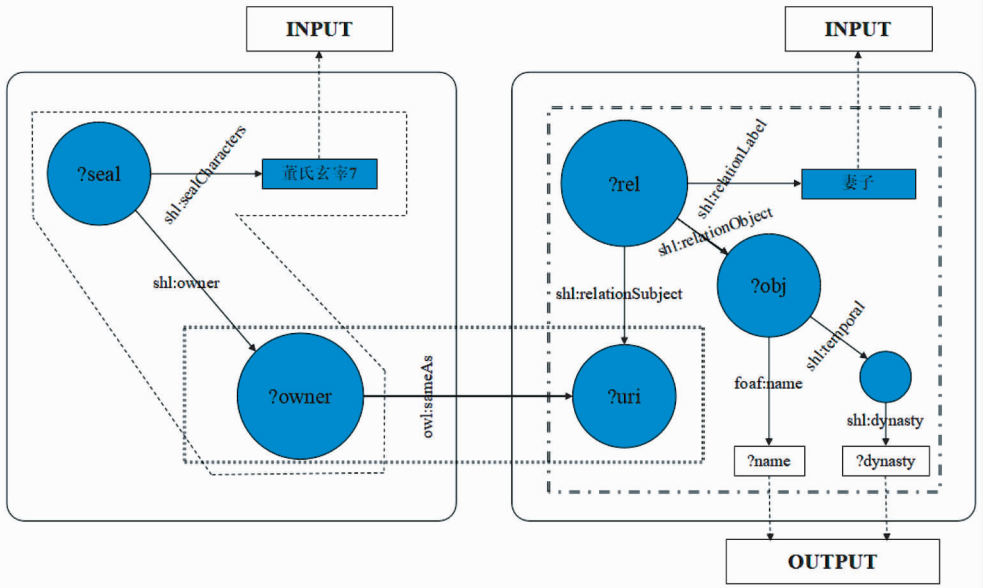


图 8 不同知识库之间隐含关系推理过程

实现图 7 推理过程的 SPARQL 联合查询语句如下:

```
PREFIX owl: <http://www.w3.org/2002/07/owl#>
PREFIX shl: <http://www.library.sh.cn/ontology/>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
SELECT ? name ? dynasty
WHERE {
  SERVICE <http://data.library.sh.cn:8890/sparql> {
    ? seal a shl:Seal;
    shl:sealCharacters "董氏玄宰7"@cht;
    shl:owner ? owner.
    ? owner owl:sameAs ? uri.
  }
```

```
SERVICE <http://cbdb.library.sh.cn/sparql> {
  ? rel a shl:Relationship;
  shl:relationLabel 妻子;
  shl:relationSubject ? uri;
  shl:relationObject ? obj.
  ? obj foaf:name ? name;
  shl:temporal/shl:dynasty ? dynasty.
  FILTER (lang(? name) = 'cht')
  FILTER (lang(? dynasty) = 'cht')
}
```

印章平台除了支持 SPARQL 联合查询,也支持基于知识图谱的可视化语义查询,通过点击知识图谱上

不同的指向,进行不同知识库之间的联合查询。如图 9 可以通过点击不同指向查询上海图书馆人名规范库、CBDB 知识库,在知识图谱上更直观地看出“董其

昌的妻子是明代成岫”,也可以通过点击不同指向查询上海图书馆古籍知识库发现董其昌相关的两部古籍《黄庭经》和《百花亭》明抄本。

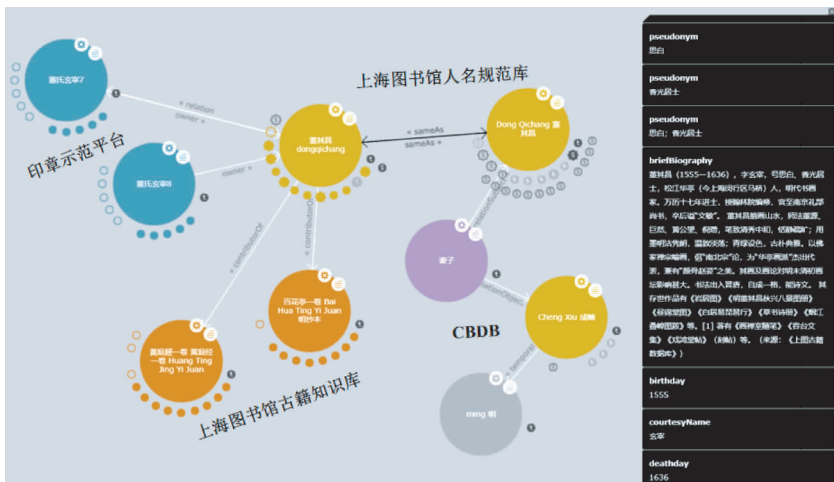


图 9 知识图谱可视化语义查询

通过印章平台获取的“董其昌的妻子是明代成岫”以及董其昌相关的两部古籍《黄庭经》和《百花亭》明抄本的发现,是传统的数据库不能实现的,这也进一步验证了图像知识组织模型对图像语义增强的可用性和有效性。

5 总结与展望

本文对国际图像互操作框架 IIIF、关联数据和知识图谱等技术在数字人文图像资源整合与知识发现领域的应用做了深入研究和探索,并围绕图像数据特征,提出了一套图像资源整合与知识发现解决方案。该解决方案从底层数据建设开始使用关联数据开放标准,并支持复用国内外开放的关联数据集,同时将 IIIF 和 OADM 两个国际通用标准相结合,实现“印章知识中心”和其他馆藏机构之间图像资源的互操作,实现了有出处来源和受控的分享与复用。模型也尝试使用 CNNs 对图像特征进行提取,实现基于机器学习的图像特征检索,同时以知识图谱的可视化方式实现多数据源的图像关联知识的语义检索与知识发现。

该解决方案可用于数字人文领域包括印章、绘画等图像资源的整合、共享与知识发现,以及数字人文研究平台和研究环境的构建,对于推动图书馆馆藏、特藏、古籍图像资源的语义化建设具有十分重要的意义。与国际知名 Linked Canvas 图像语义注释共享解决方案相比,该解决方案侧重于基于 IIIF、知识图谱(关联数据)的图像资源的整合与知识发现,在图像内容的语

义注释方面,实现了基于 OADM 数据模型的图像整体的注释,成功探索了基于 IIIF 和 OADM 数据模型的图像资源的注释、关联和共享的解决方案,可以在进一步的研究中用于在线多人协作的图像局部注释与共享,实现更深层次的知识发现,也可为科研人员提供更便捷的在线协作的学术研究环境。

参考文献:

- [1] 刘炜,叶鹰. 数字人文的技术体系与理论结构探讨[J]. 中国图书馆学报,2017,43(5):32-41.
- [2] IIIF Consortium. IIIF presentation API 2.0[EB/OL]. [2019-12-26]. <http://iiif.io/api/presentation/2.0/>.
- [3] Open Annotation Community Group. Open annotation data model. [EB/OL]. [2019-12-26]. <http://www.openannotation.org/spec/core/>.
- [4] 曾蕾,王晓光,范炜. 图档博领域的智慧数据及其在数字人文研究中的角色[J]. 中国图书馆学报,2018,44(1):17-34.
- [5] 陈涛,刘炜,单蓉蓉,等. 知识图谱在数字人文中的应用研究[J]. 中国图书馆学报,2019,45(6):1-19.
- [6] Linked canvas, engaging people, art and ideas[EB/OL]. [2019-10-26]. https://www.synaptica.com/wp-content/uploads/2015/03/Linked_Canvas_Factsheet.pdf.
- [7] ALISON A. The ‘time machine’ reconstructing ancient Venice’s social networks[J]. Nature, 2017, 7658(546):341-344.
- [8] 夏翠娟,张磊,贺晨芝. 面向知识服务的图书馆数字人文项目建设:方法、流程与技术[J]. 图书馆论坛,2018,38(1):1-9.
- [9] 曾子明,秦思琪. 面向数字人文的移动视觉搜索模型研究[J]. 情报资料工作,2018,39(6):21-28.
- [10] 侯西龙,谈国新,庄文杰,等. 基于关联数据的非物质文化遗产

知识管理研究[J]. 中国图书馆学报, 2019, 45(2): 88 - 108.

[11] 中国历代人物传记资料库 (CBDB) [EB/OL]. [2019 - 10 - 26]. <http://cbdb.library.sh.cn/>.

[12] IIIF Presentation API 1.0 [EB/OL]. [2019 - 10 - 26]. <http://iiif.io/api/search/1.0/>.

[13] TOUSCH A M, HERBIN S, AUDIBERT J Y. Semantic hierarchies for image annotation; a survey[J]. Pattern recognition, 2012, 45(1): 333 - 345.

[14] 陈涛, 刘炜, 朱庆华. 中文百科概念术语服务平台 SinoPedia 的构建研究[J]. 中国图书馆学报, 2018, 44(4): 4 - 18.

[15] 陈涛, 张永娟, 刘炜, 等. 关联数据发布的若干规范及建议[J]. 中国图书馆学报, 2019, 45(1): 34 - 46.

[16] CHEN T, ZHANG Y, WANG Z, et al. (2019) SinoPedia-A linked data services platform for decentralized knowledge base[J]. PLOS ONE, 2019, 14(8): e0219992.

[17] 印章知识中心 [EB/OL]. [2019 - 10 - 26]. <http://sinopedia.library.sh.cn:8180/seal/seal/search>.

作者贡献说明:

张永娟: 论文框架设计及论文撰写;
刘炜: 系统规划与指导, 修改建议提出;
于建荣: 论文框架指导, 修改建议提出;
陈涛: 算法指导, 系统实现。

Seal Image Resource Integration and Knowledge Discovery Based on IIIF and Semantic Knowledge Graph

Zhang Yongjuan^{1,2} Liu Wei³ Yu Jianrong² Chen Tao^{3,4}

¹ Department of Library and Information Systems, Shanghai University, Shanghai 200444

² Shanghai Institute of Nutrition and Health, Chinese Academy of Sciences / Shanghai Information Center of Life Science, Shanghai 200031

³ Shanghai Library / Shanghai Institute of Science and Technology Information, Shanghai 200031

⁴ School of Information Management, Nanjing University, Nanjing 210023

Abstract: [Purpose/significance] The image resources of digital humanities research contain a lot of information but the utilization rate is extremely low, so it cannot be effectively shared and reused in heterogeneous databases and different applications. The International Image Semantic Interoperability Framework (IIIF) breaks the barriers to image resource exchange and sharing. [Method/process] This study combined IIIF and semantic knowledge graph (linked data technology) to integrate, share and discover knowledge of image resources, reveal the relationship between resources and knowledge reasoning, and it realized semantic retrieval based on image features to assist knowledge discovery by the feature extraction and recognition of image features through CNNs algorithm. [Result/conclusion] Finally, a set of digital human image resource integration and knowledge discovery solutions was proposed, and the “Seal Knowledge Center” was constructed with the seal image resources as the application object to empirically test the feasibility and practicality of the above solutions.

Keywords: digital humanity image resource integration IIIF linked data knowledge graph knowledge discovery